

METHOD AND APPARATUS FOR UPDATING AND INVALIDATING STORE DATA

1 Cross Reference To Related Application(s)

2 This is a divisional of a copending application serial number 10/230,188, entitled
3 "METHOD AND APPARATUS FOR UPDATING AND INVALIDATING STORE
4 DATA," filed on August 29, 2002, ^{now U.S. Patent No. 6,772,316} which in turn is a divisional of application serial
5 number 09/466,306 entitled "UPDATING AND INVALIDATING STORE DATA AND
6 REMOVING STALE CACHE LINES IN A PREVALIDATED TAG CACHE
7 DESIGN," now U.S. Patent 6,470,437. These applications and the patent are hereby
8 incorporated by reference.

9 Technical Field

10 The technical field encompasses computer systems employing prevalidated cache
11 tag designs. In particular, the technical field encompasses designs to support store
12 updates and invalidates and removal of stale cache lines out of a cache.

13 Background

14 Computer systems may employ a multi-level hierarchy of memory, with relatively
15 fast, expensive but limited-capacity memory at the highest level of the hierarchy and
16 proceeding to relatively slower, lower cost but higher-capacity memory at the lowest
17 level of the hierarchy. The hierarchy may include a small fast memory called a cache,
18 either physically integrated within a processor or mounted physically close to the
19 processor for speed. The computer system may employ separate instruction caches and
20 data caches. In addition, the computer system may use multiple levels of caches. The use
21 of a cache is transparent to a computer program at the instruction level and can thus be
22 added to a computer architecture without changing the instruction set or requiring
23 modification to existing programs.

24 A cache hit occurs when a processor requests an item from a cache and the item is
25 present in the cache. A cache miss occurs when a processor requests an item from a
26 cache and the item is not present in the cache. In the event of a cache miss, the processor
27 retrieves the requested item from a lower level of the memory hierarchy. In many
28 processor designs, the time required to access an item for a cache hit is one of the primary
29 limiters for the clock rate of the processor if the designer is seeking a single cycle cache
30 access time. In other designs, the cache access time may be multiple cycles, but the
31 performance of a processor can be improved in most cases when the cache access time in

DT
2/4/05

DT
2/4/05

now, U.S. Patent No. 6,014,732 ,

1 entitled CACHE MEMORY WITH REDUCED ACCESS TIME, the disclosure of which
2 is hereby incorporated by reference.

3 The micro-architecture illustrated in Figure 3 includes a prevalidated tag cache.
4 The prevalidation imposes restrictions on how the TLBs in the micro-architecture work if
5 the computer micro-architecture is designed to maximize overall bandwidth while
6 minimizing cache load latency. The prevalidated tag cache, for example, provides very
7 fast access time for certain loads but the micro-architecture may be designed to restrict
8 the translations between virtual and physical address and restrict the distribution of
9 processing among the different cache levels. The micro-architecture may provide for fast
10 integer loads and a high bandwidth for floating point loads, for example. That is, integer
11 load data needs to have fast access timing but its working set size is generally small. To
12 optimize integer load latency, some processors provide a small but fast first level cache.
13 To provide virtual address translation and avoid address aliasing problems, some
14 processors must access the TLB to provide a physical address for checking with the cache
15 tags to determine if the data is present in the cache or not. To decrease the memory
16 latency for fast integer data access, TLB size may be limited to a small number of entries
17 (such as 16 to 32). This conflicts with the large number of entries required on processors
18 with large cache structures that could require 256 or more TLB entries.

19 In a prevalidated cache tag system, such as that shown in Figure 3, the TLB
20 entries are logically used in the cache tag to identify the cache lines. When a TLB entry
21 is removed, control is normally used to invalidate all the data in the prevalidated cache
22 tag that is associated with the removed TLB entry. However, this action may slow
23 processing since one TLB entry may map to much or all of the data cache. The TLB may
24 then be continually swapping pages in and out of memory (i.e., thrashing) instead of
25 supporting program execution.

26 Floating point data processing performance is usually limited by the memory
27 bandwidth in and out of the floating point execution units. As opposed to integer load
28 data accesses, which need a low latency, floating point accesses can usually be scheduled
29 and can therefore endure a longer latency period. Likewise, while the integer data size is
30 usually small, floating point data sets are usually very large. Ideally, TLB operations for
31 floating point load/store operations will provide both high bandwidth and large data space
32 translations (large number of TLB entries accessed). One design provides full bandwidth
33 for all memory ports and a large but slower TLB for translation of floating point requests.